

Inconsistency of Template Estimation with the Fréchet mean in Quotient Space

Loïc Devilliers* Xavier Pennec† Stéphanie Allasonnière‡

March 8, 2017

Abstract

We tackle the problem of template estimation when data have been randomly transformed under an isometric group action in the presence of noise. In order to estimate the template, one often minimizes the variance when the influence of the transformations have been removed (computation of the Fréchet mean in quotient space). The consistency bias is defined as the distance (possibly zero) between the orbit of the template and the orbit of one element which minimizes the variance. In this article we establish an asymptotic behavior of the consistency bias with respect to the noise level. This behavior is linear with respect to the noise level. As a result the inconsistency is unavoidable as soon as the noise is large enough. In practice, the template estimation with a finite sample is often done with an algorithm called max-max. We show the convergence of this algorithm to an empirical Karcher mean. Finally, our numerical experiments show that the bias observed in practice cannot be attributed to the small sample size or to a convergence problem but is indeed due to the previously studied inconsistency.

*Université Côte d'Azur, Inria, France loic.devilliers@inria.fr

†Université Côte d'Azur, Inria, France

‡Université Paris Descartes, INSERM UMRS 1138, Centre de Recherche des Cordeliers, France

Contents

1	Introduction	2
2	Inconsistency of the Template Estimation	4
3	Template estimation with the Max-Max Algorithm	5
3.1	Max-Max Algorithm Converges to a Local Minima of the Empirical Variance	5
3.2	Max-Max Algorithm is a Gradient Descent of the Variance . . .	8
4	Simulation on synthetic data	9
4.1	Max-max algorithm with a step function as template	9
4.2	Max-max algorithm with a continuous template	10
4.3	Does the max-max algorithm give us a global minimum or only a local minimum of the variance?	11
5	Discussion and Conclusion	11
A	Proof of Theorem 1	13

1 Introduction

The template estimation is a well known issue in different fields such as statistics on signals [KSW11], shape theory, computational anatomy [GMT00, JDJG04, CMT⁺04] etc. In these fields, the template (which can be viewed as the prototype of our data) can be (according to different vocabulary) shifted, transformed, wrapped or deformed due to different groups acting on data. Moreover, due to a limited precision in the measurement, the presence of noise is almost always unavoidable. These mixed effects on data lead us to study the consistency of algorithms which claim to compute the template. A popular algorithm consists in the minimization of the variance, in other words, the computation of the Fréchet mean in quotient space. This method has been already proved to be inconsistent [BC11, MHP16, DATP16]. One way to avoid the inconsistency is to use another framework, for a instance a Bayesian paradigm [CDH16]. However, if one does not want to change the paradigm, then one needs to have a better understanding of the geometrical and statistical origins of the inconsistency.

Notation: in this paper, we suppose that observations belong to a Hilbert space $(H, \langle \cdot, \cdot \rangle)$, we denote by $\|\cdot\|$ the norm associated to the dot product $\langle \cdot, \cdot \rangle$. We also consider a group of transformation G which acts isometrically on H the space of observations. This means that $x \mapsto g \cdot x$ is a linear automorphism of H , such that¹ $\|g \cdot x\| = \|x\|$, $g' \cdot (g \cdot x) = (g'g) \cdot x$ and $e \cdot x = x$ for all $x \in H$, $g, g' \in G$, where e is the identity element of G .

¹Note that in this article, $g \cdot x$ is the result of the action of g on x , and \cdot should not to be confused with the multiplication of real numbers noted \times .

The generative model is the following: we transform an unknown template $t_0 \in H$ with ϕ a random and unknown element of the group G and we add some noise $\sigma\epsilon$ with a positive noise level σ , ϵ a standardized noise: $\mathbb{E}(\epsilon) = 0$, $\mathbb{E}(\|\epsilon\|^2) = 1$. Moreover we suppose that ϵ and ϕ are independent random variables. Finally, the only observable random variable is:

$$Y = \phi \cdot t_0 + \sigma\epsilon. \quad (1)$$

If we assume that the noise is independent and identically distributed on each pixel or voxel with a standard deviation s , then $\sigma = \sqrt{N}s$, where N is the number of pixels/voxels.

Quotient space and Fréchet mean: the random transformation of the template by the group leads us to project the observation Y into the quotient space defined as the set containing all the orbit $[x] = \{g \cdot x, g \in G\}$ for $x \in H$. Because the action is isometric, the quotient space H/G is equipped with a pseudometric² defined by:

$$d_Q([x], [y]) = \inf_{g \in G} \|x - g \cdot y\| = \inf_{g \in G} \|g \cdot x - y\|.$$

The quotient pseudometric is the distance between x and y' where y' is the registration of y with respect to x . We define the variance of the random orbit $[Y]$ as the expectation of the square pseudometric between the random orbit $[Y]$ and the orbit of a point x in H :

$$F(x) = \mathbb{E}(d_Q^2([x], [Y])) = \mathbb{E}(\inf_{g \in G} \|x - g \cdot Y\|^2) = \mathbb{E}(\inf_{g \in G} \|g \cdot x - Y\|^2). \quad (2)$$

Note that $F(x)$ is well defined for all $x \in H$ because $\mathbb{E}(\|Y\|^2)$ is finite. In order to estimate the template, one often minimizes this function. If $m_\star \in H$ minimizes F , then $[m_\star]$ is called a Fréchet mean of $[Y]$. The consistency bias, noted CB , is the pseudometric between the orbit of the template $[t_0]$ and $[m_\star]$: $CB = d_Q([t_0], [m_\star])$. If such a m_\star does not exist, then the consistency bias is infinite.

Questions:

- What is the behavior of the consistency with respect to the noise?
- How to perform such a minimization of the variance? Indeed, in practice we have only a sample and not the whole distribution.

Contribution: in this article, we provide a Taylor expansion of the consistency bias when the noise level σ tends to infinity. As we do not have the whole distribution, we minimize the empirical variance given a sample. An element which minimizes the variance is called an empirical Fréchet mean. We already know that the empirical Fréchet mean converges to the Fréchet mean when the sample size tends to infinity [Zie77]. Therefore our problem is reduced to

² d_Q is called a pseudometric because $d_Q([x], [y])$ can be equal to zero even if $[x] \neq [y]$. If the orbits are closed sets then d_Q is a distance.

find an empirical Fréchet mean with a finite but sufficiently large sample. One algorithm called the max-max algorithm [AAT07] aims to compute such an empirical Fréchet mean. We establish some properties of the convergence of this algorithm. In particular, when the group is finite, the algorithm converges in a finite number of steps to an empirical Karcher mean (a local minimum of the empirical variance given a sample). This helps us to illustrate the inconsistency in this very simple framework.

Of course, generally people use a subgroup of diffeomorphisms which acts non isometrically on data such that images, landmarks etc. We believe that studying the inconsistency in this simplified framework will help us to better understand more complex situations. Moreover it is also possible to define and use isometric actions on curves [HCG⁺13, KSW11] or on surfaces [KKD⁺11] where our work can be directly applied.

This article is organized as follows: in Section 2, we study the presence of the inconsistency and we establish the asymptotic behavior when the noise parameter σ tends to ∞ . In Section 3 we detail the max-max algorithm and its properties. Finally, in Section 4 we illustrate the inconsistency with synthetic data.

2 Inconsistency of the Template Estimation

We start with the main theorem of this article which gives us an asymptotic behavior of the consistency bias when the noise level σ tends to infinity. One key notion in Theorem 1 is the concept of fixed point under the action G : a point $x \in H$ is a fixed point if for all $g \in G$, $g \cdot x = x$. We require that the support of the noise ϵ is not included in the set of fixed points. But this condition is almost always fulfilled. For instance in \mathbb{R}^n the set of fixed points under a linear group action is a null set for the Lebesgue measure (unless the action is trivial $g \cdot x = x$ for all $g \in G$ but this situation is irrelevant).

Theorem 1. *Let us suppose that the support of the noise ϵ is not included in the set of fixed points under the group action. Let Y be the observable variable defined in Equation (1). If the Fréchet mean of $[Y]$ exists, then we have the following lower and upper bounds of the consistency bias noted CB :*

$$\sigma K - 2\|t_0\| \leq CB \leq \sigma K + 2\|t_0\|, \quad (3)$$

where $K = \sup_{\|v\|=1} \mathbb{E} \left(\sup_{g \in G} \langle g \cdot v, \epsilon \rangle \right)$ is a constant which depends only on the standardised noise and on the group action. We have $K \in (0, 1]$. The consistency bias has the following asymptotic behavior when the noise level σ tends to infinity:

$$CB = \sigma K + o(\sigma) \text{ as } \sigma \rightarrow +\infty. \quad (4)$$

It follows from Equation (3) that K is the consistency bias with a null template $t_0 = 0$ and a standardised noise $\sigma = 1$. We can ensure the presence of

inconsistency as soon as the signal to noise ratio verifies $\frac{\|t_0\|}{\sigma} < \frac{K}{2}$. Moreover, if the signal to noise ratio verifies $\frac{\|t_0\|}{\sigma} < \frac{K}{3}$ then the consistency bias verifies $CB \geq \|t_0\|$. In other words, the Fréchet mean in quotient space is too far from the template: the template estimation with the Fréchet mean in quotient space is useless in this case. In [DATP16] the authors also give lower and upper bounds as a function of σ but these bounds are less informative than our current bounds. Indeed, in [DATP16] the lower bound goes to zero when the template becomes closed to fixed points. This may suggest that the consistency bias was small for this kind of template, which is not the case. The proof of Theorem 1 is postponed in Appendix A, the sketch of the proof is the following:

- $K > 0$ because the support of ϵ is not included in the set of fixed points under the action of G .
- $K \leq 1$ is the consequence of the Cauchy-Schwarz inequality.
- The proof of Inequalities (3) is based on the triangular inequalities:

$$\|m_\star\| - \|t_0\| \leq CB = \inf_{g \in G} \|t_0 - g \cdot m_\star\| \leq \|t_0\| + \|m_\star\|, \quad (5)$$

where m_\star minimizes (2): having a piece of information about the norm of m_\star is enough to deduce a piece of information about the consistency bias.

- The asymptotic Taylor expansion of the consistency bias (4) is the direct consequence of inequalities (3).

Note that Theorem 1 is absolutely not a contradiction with [KSW11] where the authors proved the consistency of the template estimation with the Fréchet mean in quotient space for all $\sigma > 0$. Indeed their noise was included in the set of constant functions which are the fixed points under their group action.

One disadvantage of Theorem 1 is that it ensures the presence of inconsistency for σ large enough but it says nothing when σ is small, in this case one can refer to [MHP16] or [DATP16].

3 Template estimation with the Max-Max Algorithm

3.1 Max-Max Algorithm Converges to a Local Minima of the Empirical Variance

Section 2 can be roughly understood as follows: if we want to estimate the template by minimising the Fréchet mean with quotient space then there is a bias. This supposes that we are able to compute such a Fréchet mean. In practice, we cannot minimise the exact variance in quotient space, because we have only a finite sample and not the whole distribution. In this section we study the estimation of the empirical Fréchet mean with the max-max algorithm. We

suppose that the group is finite. Indeed, in this case, the registration can always be found by an exhaustive search. In a compact group acting continuously, the registration also exists but is not necessarily computable without approximation. Hence, the numeric experiments which we conduct in Section 4 lead to an empirical Karcher mean in a finite number of steps.

If we have a sample: Y_1, \dots, Y_I of independent and identically distributed copies of Y , then we define the empirical variance in the quotient space:

$$F_I(x) = \frac{1}{I} \sum_{i=1}^I d_Q^2([x], [Y_i]) = \frac{1}{I} \sum_{i=1}^I \min_{g_i \in G} \|x - g_i \cdot Y_i\|^2 = \frac{1}{I} \sum_{i=1}^I \min_{g_i \in G} \|g_i \cdot x - Y_i\|^2. \quad (6)$$

The empirical variance is an approximation of the variance, indeed thanks to the law of large number we have $\lim_{I \rightarrow \infty} F_I(x) = F(x)$ for all $x \in H$. One element which minimizes globally (respectively locally) F_I is called an empirical Fréchet mean (respectively an empirical Karcher mean). For $x \in H$ and $\underline{g} \in G^I$: $\underline{g} = (g_1, \dots, g_I)$ where $g_i \in G$ for all $i \in 1..I$ we define J an auxiliary function by:

$$J(x, \underline{g}) = \frac{1}{I} \sum_{i=1}^I \|x - g_i \cdot Y_i\|^2 = \frac{1}{I} \sum_{i=1}^I \|g_i^{-1} \cdot x - Y_i\|^2.$$

The max-max algorithms iteratively minimizes the function J in the variable $x \in H$ and in the variable $\underline{g} \in G^I$:

Algorithm 1 Max-Max algorithm

Require: A starting point $m_0 \in H$, a sample Y_1, \dots, Y_I .

$n = 0$.

while Convergence is not reached **do**

Minimizing $\underline{g} \in G^I \mapsto J(m_n, \underline{g})$: we get g_i^n by registering Y_i with respect to m_n .

Minimizing $x \in H \mapsto J(x, \underline{g}^n)$: we get $m_{n+1} = \frac{1}{I} \sum_{i=1}^I g_i^n \cdot Y_i$.

$n = n + 1$.

end while

$\hat{m} = m_n$

Note that the empirical variance does not increase at each step of the algorithm since: $F_I(m_n) = J(m_n, \underline{g}^n) \geq J(m_{n+1}, \underline{g}^n) \geq J(m_{n+1}, \underline{g}^{n+1}) = F_I(m_{n+1})$. This algorithm is sensitive to the the starting point. However we remark that $m_1 = \frac{1}{I} \sum_{i=1}^I g_i \cdot Y_i$ for some $g_i \in G$, then without loss of generality, we can start from $m_1 = \frac{1}{I} \sum_{i=1}^I g_i \cdot Y_i$ for some $g_i \in G$.

Proposition 1. *As the group is finite, the convergence is reached in a finite number of steps.*

Proof. The sequence $(F_I(m_n))_{n \in \mathbb{N}}$ is non-increasing. Moreover the sequence $(m_n)_{n \in \mathbb{N}}$ takes value in a finite set which is: $\{\frac{1}{I} \sum_{i=1}^I g_i \cdot Y_i, g_i \in G\}$. Therefore, the sequence $(F_I(m_n))_{n \in \mathbb{N}}$ is stationary. Let $n \in \mathbb{N}$ such that $F_I(m_n) =$

$F_I(m_{n+1})$. Hence the empirical variance did not decrease between step n and step $n + 1$ and we have:

$$F_I(m_n) = J(m_n, \underline{g}_n) = J(m_{n+1}, \underline{g}_n) = J(m_{n+1}, \underline{g}_{n+1}) = F_I(m_{n+1}),$$

as m_n is the unique element which minimizes $m \mapsto J(m, \underline{g}_n)$ we conclude that $m_{n+1} = m_n$. \square

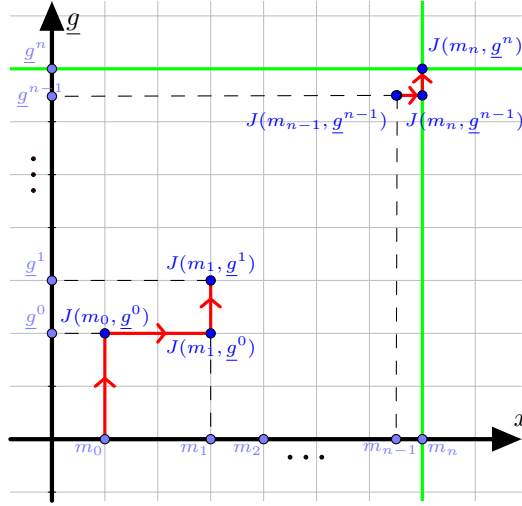


Figure 1: Iterative minimization of the function J on the two axis, the horizontal axis represents the variable in the space H , the vertical axis represents the set of all the possible registrations G^I . Once the convergence is reached, the point (m_n, g_n) is the minimum of the function J on the two axis in green. Is this point the minimum of J on its whole domain? There are two pitfalls: firstly this point could be a saddle point, it can be avoided with Proposition 2, secondly this point could be a local (but not global), this is discussed in Subsection 4.3.

This proposition gives us a shutoff parameter in the max-max algorithm: we stop the algorithm as soon as $m_n = m_{n+1}$. Let call \hat{m} the final result of the max-max algorithm. It may seem logical that \hat{m} is at least a local minimum of the empirical variance. However this intuition may be wrong: let us give a simple counterexample (but not necessarily realistic), suppose that we observe Y_1, \dots, Y_I , due to the transformation of the group it is possible that $\sum_{i=1}^n Y_i = 0$. We can start from $m_1 = 0$ in the max-max algorithm, as Y_i and 0 are already registered, the max-max algorithm does not transform Y_i . At step two, we still have $m_2 = 0$, by induction the max-max algorithm stays at 0 even if 0 is not a Fréchet or Karcher mean of $[Y]$. Because 0 is equally distant from all the points in the orbit of Y_i , 0 is called a focal point of $[Y_i]$. The notion of focal point is important for the consistency of the Fréchet mean in manifold [BP03].

Fortunately, the situation where \hat{m} is not a Karcher mean is almost always avoided due to the following statement:

Proposition 2. *Let \hat{m} be the result of the max-max algorithm. If the registration of Y_i with respect to \hat{m} is unique, in other words, if \hat{m} is not a focal point of Y_i for all $i \in 1..I$ then \hat{m} is a local minimum of F_I : $[\hat{m}]$ is an empirical Karcher mean of $[Y]$.*

Note that, if we call z the registration of y with respect to m , then the registration is unique if and only if $\langle m, z - g \cdot z \rangle \neq 0$ for all $g \in G \setminus \{e\}$. Once the max-max algorithm has reached convergence, it suffices to test this condition for \hat{m} obtained by the max-max algorithm and for Y_i for all i . This condition is in fact generic and is always obtained in practice.

Proof. We call g_i the unique element in G which register Y_i with respect to \hat{m} , for all $h \in G \setminus \{g_i\}$, $\|\hat{m} - g_i \cdot Y_i\| < \|\hat{m} - h_i \cdot Y_i\|$. By continuity of the norm we have for a close enough to m : $\|a - g_i \cdot Y_i\| < \|a - h_i \cdot Y_i\|$ for all $h_i \neq g_i$ (note that this argument requires a finite group). The registrations of Y_i with respect to m and to a are the same:

$$F_I(a) = \frac{1}{I} \sum_{i=1}^I \|a - g_i \cdot Y_i\|^2 = J(a, \underline{g}) \geq J(\hat{m}, \underline{g}) = F_I(\hat{m}),$$

because $m \mapsto J(m, \underline{g})$ has one unique local minimum \hat{m} . □

3.2 Max-Max Algorithm is a Gradient Descent of the Variance

In this Subsection, we see that the max-max algorithm is in fact a gradient descent. The gradient descent is a general method to find the minimum of a differentiable function. Here we are interested in the minimum of the variance F : let $m_0 \in H$ and we define by induction the gradient descent of the variance $m_{n+1} = m_n - \rho \nabla F(m_n)$, where $\rho > 0$ and F the variance in the quotient space. In [DATP16] the gradient of the variance in quotient space for m a regular point was computed (m is regular as soon as $g \cdot m = m$ implies $g = e$), this leads to:

$$m_{n+1} = m_n - 2\rho [m_n - \mathbb{E}(g(Y, m_n) \cdot Y)],$$

where $g(Y, m_n)$ is the almost-surely unique element of the group which register Y with respect to m_n . Now if we have a set of data Y_1, \dots, Y_n we can approximate the expectation which leads to the following approximated gradient descent:

$$m_{n+1} = m_n(1 - 2\rho) + \rho \frac{2}{I} \sum_{i=1}^I g(Y_i, m_n) \cdot Y_i,$$

now by taking $\rho = \frac{1}{2}$ we get $m_{n+1} = \frac{1}{I} \sum_{i=1}^I g(Y_i, m_n) \cdot Y_i$. So the approximated gradient descent with $\rho = \frac{1}{2}$ is exactly the max-max algorithm. But the max-max algorithm is proven to be converging in a finite number of steps which is not the case for gradient descent in general.

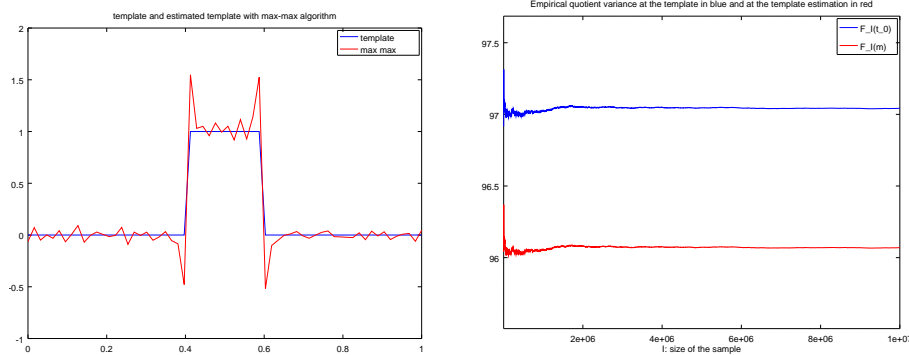
4 Simulation on synthetic data

In this Section³, we consider data in an Euclidean space \mathbb{R}^N equipped with its canonical dot product $\langle \cdot, \cdot \rangle$, and $G = \mathbb{Z}/N\mathbb{Z}$ acts on \mathbb{R}^N by circular permutation on coordinates:

$$(\bar{k} \in \mathbb{Z}/N\mathbb{Z}, (x_1, \dots, x_N) \in \mathbb{R}^N) \mapsto (x_{1+k}, x_{2+k}, \dots, x_{N+k}),$$

where indexes are taken modulo N . This space models the discretization of functions with N points. This action is found in [AAT07] and used for neuro-electric signals in [HCG⁺13]. The registration between two vectors can be made by an exhaustive research but it is faster with the fast Fourier transform [CT65].

4.1 Max-max algorithm with a step function as template



(a) Example of a template (a step function) and the template estimation with a sample size 10^5 in \mathbb{R}^{64} , ϵ is Gaussian noise and $\sigma = 10$. At the discontinuity points of the template, we observe a Gibbs-like phenomena.

(b) Variation of $F_I(t_0)$ (in blue) and of $F_I(\hat{m})$ (in red) as a function of I the size of the sample. Since convergence is already reached, $F(\hat{m})$, which is the limit of red curve, is below $F(t_0)$: $F(t_0)$ is the limit of the blue curve. Due to the inconsistency, \hat{m} is an example of point such that $F(\hat{m}) < F(t_0)$.

Figure 2: Template t_0 and template estimation \hat{m} on Fig. 2(a). Empirical variance at the template and the template estimation with the max-max algorithm as a function of the size of the sample on Fig. 2(b).

We display an example of a template and the template estimation with the max-max algorithm on Fig 2(a). Note that this experiment was already conducted in [AAT07]. But no explanation of the appearance of the bias was provided. On the opposite, we know from the precedent Section that the max-max result is an empirical Karcher mean, and that this result can be obtained in a finite number of steps. Taking $\sigma = 10$ may seem extremely high, however

³The code used in this Section is available at <http://loic.devilliers.free.fr/ipmi.html>.

the standard deviation of the noise at each point is not 10 but $\frac{\sigma}{\sqrt{N}} = 1.25$ which is not so high.

The sample size is 10^5 , and the algorithm stopped after 94 steps, and \hat{m} the estimated template (in red on the Fig. 2(a)) is not a focal points of the orbits $[Y_i]$, then Proposition 2 applies. We call empirical bias (noted EB) the quotient distance between the true template and the point \hat{m} given by the max-max result. On this experiment we have $\frac{EB}{\sigma} \simeq 0.11$. Of course, one could think that we estimate the template with an empirical bias due to a too small sample size which induces fluctuation. To reply to this objection, we keep in memory \hat{m} obtained with the max-max algorithm. If there was no inconsistency then we would have $F(t_0) \leq F(\hat{m})$. We do not know the value of the variance F at these points, but thanks to the law of large number, we know that:

$$F(t_0) = \lim_{I \rightarrow \infty} F_I(t_0) \text{ and } F(\hat{m}) = \lim_{I \rightarrow \infty} F_I(\hat{m}),$$

Given a sample, we compute $F_I(t_0)$ and $F_I(\hat{m})$ thanks to the definition of the empirical variance F_I (6). We display the result on Fig. 2(b), this tends to confirm that $F(t_0) > F(\hat{m})$. In other words, the variance at the template is bigger than the variance at the point given by the max-max algorithm.

4.2 Max-max algorithm with a continuous template

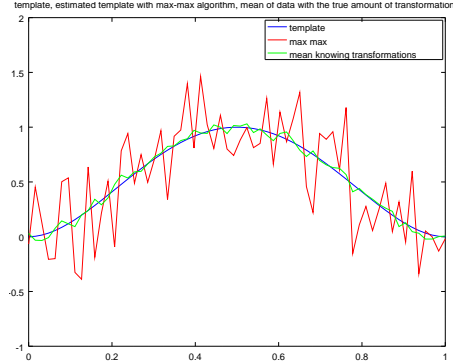


Figure 3: Example of an other template (here a discretization of a continuous function) and the template estimation with a sample size 10^3 in \mathbb{R}^{64} (in red), ϵ is Gaussian noise and $\sigma = 10$. Even with a continuous function the inconsistency appears. In green we compute the mean of data with the true amount of transformations.

Figure 2(a) shows that the main source of the inconsistency was the discontinuity of the template. We could think that a continuous template leads to consistency. But it is not the case, even with a large number of observations created from a continuous template we do not observe a convergence to the template see Fig. 3, the empirical bias satisfies $\frac{EB}{\sigma} = 0.25$. If we knew the

original transformations we could invert the transformations on data and take the mean, that is what we deed in green on Fig. 3. We see that with a sample size 10^3 , the mean gives us almost the good result since we have in that case $\frac{EB}{\sigma} = 0.03$.

4.3 Does the max-max algorithm give us a global minimum or only a local minimum of the variance?

Proposition 2 tells us that the output of the max-max algorithm is a Karcher mean of the variance, but we do not know that if it is Fréchet mean of the variance. In other words, is the output a global minimum of the variance? In fact, F_I has a lot of local minima which are not global. Indeed we can use the max-max algorithm with different starting points and we observe different outputs (which are all local minima thanks to Proposition 2) with different empirical variance (result non shown).

5 Discussion and Conclusion

We provided an asymptotic behavior of the consistency bias when the noise level σ tends to infinity, as a consequence, the inconsistency cannot be neglected when σ is large. However we have not answered this question: can the inconsistency be neglected? When the noise level is small enough, then the consistency bias is small [MHP16, DATP16], hence it can be neglected. Note that the quotient space is not a manifold, this prevents us to use *a priori* the Central Limit theorem for manifold proved in [BP03]. But if the Central Limit theorem could be applied to quotient space, the fluctuations induce an error which would be approximately equal to $\frac{\sigma}{\sqrt{I}}$ and if $K \ll \frac{1}{\sqrt{I}}$, then the inconsistency could be neglected because it is small compared to fluctuation.

If the Hilbert Space is a functional space, for instance $L^2([0, 1])$, in practice, we never observe the whole function, only a finite number values of this function. One can model these observable values on a grid. When the resolution of the grid goes to zero, one can show the consistency [PZ16] by using the Fréchet mean with the Wasserstein distance on the space of measures rather than in the space of functions. But in (medical) images the number of pixels or voxels is finite.

Finally, in a future work one needs to study the template estimation with non isometric action. But we can already learn from this work: in the numerical experiments we led, we have seen that the template estimated is more detailed that the true template. The intuition is that the estimated template in computational anatomy with a group of diffeomorphisms is also more detailed. But the true template is almost always unknown. It is then possible that one think that the computation of the template succeeded to capture small details of the template while it is just an artifact due to the inconsistency. Moreover in order to tackle this question, one needs to have a good modelisation of the noise, for

instance in [KSW11], the observations are curves, what is a relevant noise in the space of curves?

References

- [AAT07] Stéphanie Allasonnière, Yali Amit, and Alain Trouvé. Towards a coherent statistical framework for dense deformable template estimation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 69(1):3–29, 2007.
- [BC11] Jérémie Bigot and Benjamin Charlier. On the consistency of fréchet means in deformable models for curve and image analysis. *Electronic Journal of Statistics*, 5:1054–1089, 2011.
- [BP03] Rabi Bhattacharya and Vic Patrangenaru. Large sample theory of intrinsic and extrinsic sample means on manifolds. i. *Annals of statistics*, pages 1–29, 2003.
- [CDH16] Wen Cheng, Ian L Dryden, and Xianzheng Huang. Bayesian registration of functions and curves. *Bayesian Analysis*, 11(2):447–475, 2016.
- [CMT⁺04] Timothy F Cootes, Stephen Marsland, Carole J Twining, Kate Smith, and Christopher J Taylor. Groupwise diffeomorphic non-rigid registration for automatic model building. In *European conference on computer vision*, pages 316–327. Springer, 2004.
- [CT65] James W Cooley and John W Tukey. An algorithm for the machine calculation of complex fourier series. *Mathematics of computation*, 19(90):297–301, 1965.
- [DATP16] L. Devilliers, S. Allasonnière, A. Trouvé, and X. Pennec. Template estimation in computational anatomy: Fréchet means in top and quotient spaces are not consistent. *ArXiv e-prints*, August 2016.
- [GMT00] A. Guimond, J. Meunier, and J.-P. Thirion. Average brain models: A convergence study. *Computer Vision and Image Understanding*, 77(2):192–210, 2000.
- [HCG⁺13] Sebastian Hitziger, Maureen Clerc, Alexandre Gramfort, Sandrine Saillet, Christian Bénar, and Théodore Papadopoulos. Jitter-adaptive dictionary learning-application to multi-trial neuroelectric signals. *arXiv preprint arXiv:1301.3611*, 2013.
- [JDJG04] Sarang Joshi, Brad Davis, Matthieu Jomier, and Guido Gerig. Unbiased diffeomorphic atlas construction for computational anatomy. *NeuroImage*, 23:S151–S160, 2004.

- [KKD⁺11] Sebastian Kurtek, Eric Klassen, Zhaohua Ding, Malcolm J Avison, and Anuj Srivastava. Parameterization-invariant shape statistics and probabilistic classification of anatomical surfaces. In *Biennial International Conference on Information Processing in Medical Imaging*, pages 147–158. Springer, 2011.
- [KSW11] Sebastian A Kurtek, Anuj Srivastava, and Wei Wu. Signal estimation under random time-warpings and nonlinear signal alignment. In *Advances in Neural Information Processing Systems*, pages 675–683, 2011.
- [MHP16] Nina Miolane, Susan Holmes, and Xavier Pennec. Template shape estimation: correcting an asymptotic bias. *arXiv preprint arXiv:1610.01502*, 2016.
- [PZ16] Victor M Panaretos and Yoav Zemel. Amplitude and phase variation of point processes. *The Annals of Statistics*, 44(2):771–812, 2016.
- [Zie77] Herbert Ziezold. On expected figures and a strong law of large numbers for random elements in quasi-metric spaces. In *Transactions of the Seventh Prague Conference on Information Theory, Statistical Decision Functions, Random Processes and of the 1974 European Meeting of Statisticians*, pages 591–602. Springer, 1977.

A Proof of Theorem 1

Proof. In the proof, we note by S the unit sphere in H . In order to prove that $K > 0$, we take x in the support of ϵ such that x is not a fixed point under the action of G . It exists $g_0 \in G$ such that $g_0 \cdot x \neq x$. We note $v_0 = \frac{g_0 \cdot x}{\|x\|} \in S$, we have $\langle v_0, g_0 \cdot x \rangle = \|x\| > \langle v_0, x \rangle$ and by continuity of the dot product it exists $r > 0$ such that: $\forall y \in B(x, r) \quad \langle v_0, g_0 \cdot y \rangle > \langle v_0, y \rangle$ as x is in the support of ϵ we have $\mathbb{P}(\epsilon \in B(x, r)) > 0$, it follows:

$$\mathbb{P} \left(\sup_{g \in G} \langle v_0, g \cdot \epsilon \rangle > \langle v_0, \epsilon \rangle \right) > 0. \quad (7)$$

Thanks to Inequality (7) and the fact that $\sup_{g \in G} \langle v_0, g \cdot \epsilon \rangle \geq \langle v_0, \epsilon \rangle$ we have:

$$K = \sup_{v \in S} \mathbb{E} \left(\sup_{g \in G} \langle v, g \cdot \epsilon \rangle \right) \geq \mathbb{E} \left(\sup_{g \in G} \langle v_0, g \cdot \epsilon \rangle \right) > \mathbb{E}(\langle v_0, \epsilon \rangle) = \langle v_0, \mathbb{E}(\epsilon) \rangle = 0.$$

Using the Cauchy-Schwarz inequality: $K \leq \sup_{v \in S} \mathbb{E}(\|v\| \times \|\epsilon\|) \leq \mathbb{E}(\|\epsilon\|^2)^{\frac{1}{2}} = 1$. We now prove Inequalities (3). The variance at λv for $v \in S$ and $\lambda \geq 0$ is:

$$F(\lambda v) = \mathbb{E} \left(\inf_{g \in G} \|\lambda v - g \cdot Y\|^2 \right) = \lambda^2 - 2\lambda \mathbb{E} \left(\sup_{g \in G} \langle v, g \cdot Y \rangle \right) + \mathbb{E}(\|Y\|^2). \quad (8)$$

Indeed $\|g \cdot Y\| = \|Y\|$ thanks to the isometric action. We note $x^+ = \max(x, 0)$ the positive part of x and $h(v) = \mathbb{E}(\sup_{g \in G} \langle v, g \cdot Y \rangle)$. The $\lambda \geq 0$ which⁴ minimizes (8) is $h(v)^+$ and the minimum value of the variance restricted to the half line \mathbb{R}^+v is $F(h(v)^+v) = \mathbb{E}(\|Y\|^2) - (h(v)^+)^2$. To find $[m_\star]$ the Fréchet mean of $[Y]$, we need to maximize $(h(v)^+)^2$ with respect to $v \in S$: $m_\star = h(v_\star)v_\star$ with⁵ $v_\star \in \operatorname{argmax}_{v \in S} h(v)$. As we said in the sketch of the proof we are interested in getting a piece of information about the norm of $\|m_\star\|$ we have: $\|m_\star\| = h(v_\star) = \sup_{v \in S} h$. Let $v \in S$, we have: $-\|t_0\| \leq \langle v, g\phi \cdot t_0 \rangle \leq \|t_0\|$ because the action is isometric. Now we decompose $Y = \phi \cdot t_0 + \sigma\epsilon$ and we get:

$$\begin{aligned} h(v) &= \mathbb{E} \left(\sup_{g \in G} \langle v, g \cdot Y \rangle \right) = \mathbb{E} \left(\sup_{g \in G} (\langle v, g \cdot \sigma\epsilon \rangle + \langle v, g\phi \cdot t_0 \rangle) \right) \\ h(v) &\leq \mathbb{E} \left(\sup_{g \in G} (\langle v, g \cdot \sigma\epsilon \rangle + \|t_0\|) \right) = \sigma \mathbb{E} \left(\sup_{g \in G} \langle v, g \cdot \epsilon \rangle \right) + \|t_0\| \\ h(v) &\geq \mathbb{E} \left(\sup_{g \in G} (\langle v, g \cdot \sigma\epsilon \rangle) - \|t_0\| \right) = \sigma \mathbb{E} \left(\sup_{g \in G} \langle v, g \cdot \epsilon \rangle \right) - \|t_0\|. \end{aligned}$$

By taking the biggest value in these inequalities with respect to $v \in S$, by definition of K we get:

$$-\|t_0\| + \sigma K \leq \|m_\star\| \leq \|t_0\| + \sigma K. \quad (9)$$

Thanks to (9) and to (5), Inequalities (3) are proved. \square

⁴Indeed we know that $x \in \mathbb{R}^+ \mapsto x^2 - 2bx + c$ reaches its minimum at the point $x = b^+$ and $f(b^+) = c - (b^+)^2$.

⁵Note that we remove the positive part and the square because $\operatorname{argmax} h = \operatorname{argmax} (h^+)^2$ since h takes a non negative value (indeed $h(v) \geq \mathbb{E}(\langle v, \phi \cdot t_0 + \epsilon \rangle) = \langle v, \mathbb{E}(\phi \cdot t_0) \rangle$ and this last quantity is non negative for at least one $v \in S$).